# Contigi csb10a\_v1:

pliki ze strony http://csgenome.sggw.pl/en-us/resources/

Fragmenty plików:

## - Gene\_Annotation\_v1.xlsx:

	Α	В	С	D	Е	F	G	Н	T I	J	K
1	Lp	Gene	Start_Gene	Stop_Gene	Strand	ACC_NO	Species	Protein Product	goBP	goMF	goCC
2	3	gene_1#CSB10A_v1_contig_3	1043	5674				histone acetyltransferase	0	nucleic acid binding	0
3	4	gene_1#CSB10A_v1_contig_4	3352	16648	+	XP_002323740	Populus trichocarpa	histone acetyltransferase	0	nucleic acid binding	0
4	9	gene_2#CSB10A_v1_contig_5	13220	partial	+	XP_002523117	Ricinus communis	galacturan 1,4-alpha-galacturonidase, pu	carbohydrate r	polygalacturonase ad	0
5	12	gene_1#CSB10A_v1_contig_8	9819	14148	+	BAH57245	Arabidopsis thaliana		0	0	0
6	15	gene_4#CSB10A_v1_contig_8	35574	33792	-	BAF07113	Oryza sativa Japonica	a Group Spermatophyta; Magnoliophyta; Li	0	0	0
7	18	gene_2#CSB10A_v1_contig_10	9561	13870	+	NP_193037		oxidoreductase, zinc-binding dehydrogen	metabolic prod	catalytic activity  bindi	vacuole  pl:
8	19	gene_3#CSB10A_v1_contig_10	23201	24610	+	ABA81857	Solanum tuberosum	ripening regulated protein-like	metabolic prod	catalytic activity  bindi	0
9	20	gene_4#CSB10A_v1_contig_10	26895	25065	-	AAF82612	Prunus dulcis	self-incompatibility associated ribonuclea	0	RNA binding  ribonud	0
10	23	gene_2#CSB10A_v1_contig_12	2302	9209	+	CAB95025	Nicotiana tabacum	pectin methylesterase	cell wall modif	enzyme inhibitor activ	cell wall
11	26	gene 4#CSR104 v1 contig 12	20046	20702	+	XP 002527622	Picinus communis	Anthoryanin 5-aromatic acultraneferace	0	acultraneforace activi	. 0

#### - csb10a v1 predictions.fasta:

>gene\_2#CSB10A\_v1\_contig\_10146 MRRHFRDWWLGAGFDVPNYFSPIFFQIGNSATGFFIIFALIAAVA... >gene\_1#CSB10A\_v1\_contig\_10147 MAIEIQLINLGMQWPELLKLLSSSFGAERNSALLTYFFLSTMPELYT

# - csb10a\_v1\_upstream\_1000:

>gene\_1#CSB10A\_v1\_contig\_3
AAGGACCTTCTCATTCCAATTCTTTTGAATATAAGCCCTTGAGCTTTTCCCTAAAC...
>gene\_1#CSB10A\_v1\_contig\_4
GACTGATTATATATTGGTCCGTTCTGCAAAGGGGAAGCTTTCCCTGAGACGTGTT...

### -csb10a\_v1\_contig.fsa:

>CSB10A\_v1\_contig\_1

 ${\tt AAGATGAGGTGTTTATAGGACTTTAGTTTTCAAAAATTGAAATTTTAAAATGAGAT...}$ 

>CSB10A v1 contig 4

GTTCAATCGGTTGCTTGGCTAACGGAAAAACAAGAGGTCGGACGAAGA...

#### - GFF\_to\_Predictions\_ID:

ContigID	Model	Feature	Start	Stop	Strand	-	GFF_Genel	D Predictions_GeneID
CSB10A_v1_contig_1	GeneMark.hmm	stop_codor	n 11744	11746	-	0	1_g	gene_1#CSB10A_v1_contig_1
CSB10A_v1_contig_1	GeneMark.hmm	CDS	11744	11760	-	1	1_g	gene_1#CSB10A_v1_contig_1
CSB10A_v1_contig_1	GeneMark.hmm	CDS	11924	11942	-	0	1_g	gene_1#CSB10A_v1_contig_1
CSB10A_v1_contig_1	GeneMark.hmm	start_codo	n 11940	11942	-	0	1_g	gene_1#CSB10A_v1_contig_1

- 1. posiadamy sekwencje contigów (plik csb10a\_v1\_contig.fsa) => posiadamy długości contigów
- 2. znamy położenie genów względem contigów (kolumny Start\_Gene i Stop\_Gene w pliku Gene\_Annotation\_v1.xlsx)

ponieważ pkt1 i pkt2 => znamy również sekwencje genów

- 3. posiadamy drugie źródło sekwencji genów (plik csb10a\_v1\_upstream\_1000.fasta), w którym jest więcej genów niż w danych z pkt2. Dla dodatkowych genów nie mamy pozycji Start\_Gene i Stop Gene.
- 4. Jak interpretować słowo 'partial' w niektórych kolumnach (Start\_Gene i Stop\_Gene w pliku Gene Annotation v1.xlsx) ?
- 5. Jak interpretować dane gdy Start\_Gene > Stop\_Gene (np. gene\_2#CSB10A\_v1\_contig\_10146 Gene\_Annotation\_v1.xlsx) ?

- 6. Jakie informacje zawiera plik csb10a\_v1\_predictions.fasta?
- 7. Nie wiemy na jakich chromosomach i gdzie leżą contigi.
- 8. Nie znamy długości chromosomów.

# Contigi z markerami BAM:

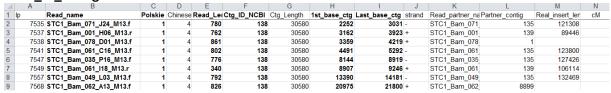
pliki https://drive.google.com/drive/folders/0B1KskZHvmCFpTEYtYlJVNnJYWk0?usp=sharing

#### Fragmenty plików:

#### - links3

	Α	В	С	D		E	F	F G		Н	1	J	K	L	M	N
1	ncbi_ctg_id	lp_oryg	#super_id	num_bases_	in_super	num_contigs_in_super	ordinal	_contig_id	length	_of_contig	estimated_gap_bet	f estimated_gap	if	start	stop	j
2	1	1	C	1	2496462	. 46	3 1		0	23626	0	-299		0	23626	2496462
3	2	2	C	1	2496462	. 46	3 2	2	1	1500	-299	56167		23327	24827	
4	3	3	C	1	2496462	! 46	3 3	3	2	8044	56167	-6380	56167	80994	89038	
5	4	4	C	1	2496462	! 46	3 4	1	3	47152	-6380	-3469		82658	129810	
6	5	5	C	1	2496462	! 46	6 5	5	4	14055	-3469	-300		126341	140396	
7	6	6	C	1	2496462	! 46	6 6	3	5	1500	-300	51167		140096	141596	
8	7	7	C	1	2496462	! 46	3 7	7	6	1500	51167	-299	51167	192763	194263	
9	8	8	C	1	2496462	! 46	3 8	3	7	36453	-299	-14784		193964	230417	
10	9	9	C		2496462	. 46	3 9	)	8	20179	-14784	7822		215633	235812	

## - STC\_vs\_contigs1



- 1. Wiemy który contig leży na którym chromosomie (Ctg\_ID\_NCBI i Polskie chromosomy w pliku STC\_vs\_contigs1).
- 2. Nie znamy długości chromosomów.
- 3. Nie wiemy gdzie contigi leżą na chromosomach pozycje start i stop (plik links3) są bezwzględne, nie dzielą organizmu na chromosomy.
- 4. Jak powiązać Ctg\_ID\_NCBI z contigami csb10a\_v1?
- 5. Nie znamy sekwencji contigów identyfikowanych Ctg ID NCBI
- 6. Nie znamy sekwencji BAMów
- 7. Znamy położenie BAMów względem contigów (o id Ctg\_ID\_NCBI) od 1st\_base\_ctg do Last\_base\_ctg (plik STC\_vs\_contigs1)